

Problem and Assumptions

Regularized Optimization

 $\min_{x \in \mathbb{R}^n} F(x) = f(x) + R(x)$ (1)

• $f: \mathbf{M}$ -smooth & μ -strongly convex convex:

$$f(x+h) \le f(x) + \langle \nabla f(x), h \rangle + \frac{1}{2} \langle \mathbf{M}h, h \rangle$$

 $f(x) + \langle \nabla f(x), h \rangle + \frac{\mu}{2} \|h\|^2 \le f(x+h)$

(natural assumptions for ERM with linear predictors)

• R non-smooth, convex & proximable

New Oracle: Gradient Sketch

We do not have direct access to $\nabla f(x)$. Instead, we have access to a random linear transformation of the gradient:

$$\mathbf{S}^{\top} \nabla f(x) \in \mathbb{R}^{b}, \qquad \mathbf{S} \sim \mathcal{D}$$
 (2)

• S: random $n \times b$ matrix (b small)

• \mathcal{D} : distribution from which **S** is drawn

Goal

Design a proximal stochastic gradient-type method for solving (1) using the gradient sketch oracle (2).

Simple Algorithmic Idea

$$x^{k+1} = \operatorname{prox}_{\alpha R}(x^k - \alpha g^k), \qquad (3)$$

$$\alpha = \text{stepsize}; \ g^k = \text{a "nice" estimator of } \nabla f(x^k).$$

How to design a good gradient estimator q^k ?

Key Challenges:

• In the case when \mathcal{D} is a distribution over standard basis vectors e_1, \ldots, e_n in \mathbb{R}^n , i.e., if we have access to random partial derivatives of f, then we can use

$$g^k = e_i^\top \nabla f(x^k) e_i,$$

and (3) reduces to proximal randomized coordinate descent (CD). However, CD does not work with non-separable regularizers R. So, we have an issue even in this simple case! How to resolve it?

• How to deal with gradient sketches coming from any distribution \mathcal{D} ?

Resolution: The SEGA estimator. We will iteratively learn an unbiased variance-reduced estimator g^k of the gradient $\nabla f(x^k)$ by incorporating the latest information provided by the gradient sketch.

SEGA: Variance Reduction via Gradient Sketching

Filip Hanzely¹

Konstantin Mishchenko¹

¹KAUST

²University of Edinburgh

³Moscow Institute of Physics and Technology

Constructing the SEGA Estimator

SEGA Estimator **1** Ask oracle for a gradient sketch at x^k : $\mathbf{S}_k^\top \nabla f(x^k)$ **2** Define h^{k+1} as the closest (in some energy norm $||h||_{\mathbf{B}}^2 \stackrel{\text{def}}{=} h^{\top} \mathbf{B} h$, where $\mathbf{B} \succ 0$) vector to h^k consistent with the gradient sketch: $h^{k+1} = \arg\min_{h \in \mathbb{R}^n} \|h - h^k\|_{\mathbf{B}}^2$ subject to $\mathbf{S}_k^{\mathsf{T}} h = \mathbf{S}_k^{\mathsf{T}} \nabla f(x^k)$ Closed-form solution of (4): $h^{k+1} = h^k + \mathbf{B}^{-1} \mathbf{Z}_k (\nabla f(x^k) - h^k); \qquad \mathbf{Z}_k \stackrel{\text{def}}{=} \mathbf{S}_k \left(\mathbf{S}_k^\top \mathbf{S}_k \right)^{\dagger} \mathbf{S}_k^\top$ **3** Define the **SEGA** estimator: $g^{k} = h^{k} + \theta_{k} \mathbf{B}^{-1} \mathbf{Z}_{k} (\nabla f(x^{k}) - h^{k})$ (5) $(\theta_k \text{ is a random variable ensuring that } g^k \text{ is unbiased})$

Key property: As $x_k \to x^*$, we get $g^k \to 0$, and hence **SEGA** estimator is variance-reduced.

Variants:

• biasSEGA estimator: use h^{k+1} instead of g^k

• subspaceSEGA estimator: If $f(x) = \phi(\mathbf{A}x)$ for some matrix $\mathbf{A} \in \mathbb{R}^{d \times n}$, we can improve the SEGA estimator by exploiting the fact that ∇f lies in Range(\mathbf{A}^{\top}). We do this by adding the constraint $h \in \text{Range}(\mathbf{A}^{+})$ to (4).

SEGA (SkEtched GrAdient descent)

SEGA = Method (3) + SEGA estimator (5)biasSEGA = Method (3) + biasSEGA estimator (4)subspaceSEGA = Method (3) + subspaceSEGA estimator

Convergence of SEGA

Let \mathcal{D} be the uniform distribution over standard basis vectors $e_1, \ldots, e_n \in \mathbb{R}^n$, and choose $\mathbf{B} = \mathbf{I}$. Then with stepsize $\alpha = \Omega(\frac{1}{n\lambda_{m}})$ and some constant $\sigma > 0$, we have

$$\mathbb{E}\left[\Phi^{\kappa}\right] \leq (1 - \alpha \mu)^{\kappa} \mathbb{E}\left[\Phi^{0}\right],$$

where $\Phi^{k} \stackrel{\text{def}}{=} \|x^{k} - x^{*}\|^{2} + \sigma \alpha \|h^{k} - \nabla f(x^{*})\|^{2}, x^{*} =$ $\operatorname{arg\,min}_x F(x).$

- Note that $x^k \to x^*$ and $h^k \to \nabla f(x^*)$
- General convergence result for any $\mathbf{B} \succ \mathbf{0}$ and any \mathcal{D} can be found in the paper [1].
- subspaceSEGA: If \mathcal{D} samples from the columns of \mathbf{A}^{+} , the rate can be $\Omega(\frac{n}{d})$ faster than standard SEGA.
- For coordinate sketches, we designed an accelerated SEGA, and established accelerated rate (read next).





Bottom plot: R is the indicator function of the unit ball. While CD does not converge, SEGA does!

Peter Richtárik^{1, 2, 3}

Experiments

1. SEGA vs Random Direct Search (RDS) [2] (coordinate and Gaussian sketches) for derivative-free optimization







Setup:

- S are column submatrices of the identity matrix
- Probability vector $p \in \mathbb{R}^n$: $p_i \stackrel{\text{def}}{=} \operatorname{Prob}(e_i \in \mathbf{S})$
- Probability matrix $\mathbf{P} \in \mathbb{R}^{n \times n}$: $\mathbf{P}_{ij} \stackrel{\text{def}}{=} \operatorname{Prob}(e_i \in \mathbf{S}, e_j \in \mathbf{S})$ • ESO vector $v \in \mathbb{R}^n$ (for mini-batching) defined by:

Algori

- 2: **for**

7: **en**

Up to the constant factors 8.55 and 9.5, these rates are exactly the same as the rates of coordinate descent [3] and accelerated coordinate descent [4, 5]. So, we extend the reach of coordinate descent methods to problem (1) with a non-separable regularizer (e.g., arbitrary convex constraint)

- [1] Filip Hanzely, Konstantin Mishchenko, and Peter Richtárik. SEGA: Variance reduction via gradient sketching. In NeurIPS, 2018.
- [2] El Houcine Bergou, Peter Richtárik, and Eduard Gorbunov. Stochastic three point method for minimizing nonconvex, convex and strongly convex functions. Manuscript, 2018.
- [4] Zeyuan Allen-Zhu, Zheng Qu, Peter Richtárik, and Yang Yuan. Even faster accelerated coordinate descent using non-uniform sampling. In *ICML*, 2016.
- [5] Filip Hanzely and Peter Richtárik. Accelerated coordinate descent with arbitrary sampling and best rates for minibatches. arXiv:1809.09354, 2018.





SEGA with Coordinate Sketches

$$\mathbf{P} \bullet \mathbf{M} \preceq \operatorname{Diag}(p \bullet v)$$

Acceleration: For coordinate sketches we also designed an accelerated variant of SEGA:

thm Accelerated SEGA (ASEGA)

$$= y^{0} = z^{0} \in \mathbb{R}^{n}; h^{0} \in \mathbb{R}^{n}; \text{ params } \alpha, \beta, \tau, \mu > 0$$

$$k = 1, 2, \dots \text{ do}$$

$$x^{k} = (1 - \tau)y^{k-1} + \tau z^{k-1}$$
Sample $\mathbf{S}_{k} \sim \mathcal{D}$, and compute g^{k} and h^{k+1}

$$y^{k} = x^{k} - \alpha p^{-1} \bullet g^{k}$$

$$z^{k} = \frac{1}{1+\beta\mu}(z^{k} + \beta\mu x^{k} - \beta g^{k})$$
d for

Rates: We prove the following iteration complexity bounds of SEGA and ASEGA with coordinate sketches:

Method	Complexity
SEGA	$\frac{1}{855}$, $\operatorname{Tr}(\mathbf{M})$ $\frac{1}{100}$
importance sampling	$\mu \log \frac{108}{\epsilon}$
SEGA	8 55 $\cdot \left(\max \cdot \frac{v_i}{v_i} \right) \log \frac{1}{v_i}$
arbitrary sampling	$(\max_{i \in p_i \mu}) \log \epsilon$
ASEGA	$0.8 \cdot \frac{\sum_i \sqrt{\mathbf{M}_{ii}}}{\log 1} \log \frac{1}{2}$
importance sampling	$\sqrt{\mu}$ $\log \epsilon$
ASEGA	$0.8.$ $\sqrt{\max \frac{v_i}{v_i}} \log \frac{1}{2}$
arbitrary sampling	$\int \frac{3.6}{\sqrt{11}} \sqrt{\frac{11}{p_i^2 \mu}} \log \frac{1}{\epsilon}$

References

- [3] Peter Richtárik and Martin Takáč.
- On optimal probabilities in stochastic coordinate descent methods. *Optimization Letters*, 10(6):1233–1243, 2016.